

Reinforcement Learning

Consider a dynamic system (not static) and its evolution over time.

We have the concept of **states**: a state is used to represent the evolution over time of the dynamic system. We also have other elements: **observations**, **state transition model** (function that describes how the system evolve over time), **noise** (used to gather everything that is not explicitly modeled).

Reasoning: given the model (f, h) and the current state x_k , predict the future (x_{k+T}, z_{k+T}) assuming we have all the necessary information about the dynamic system; if we know the transition function that tells us what the next state will be given the current one and an action to make, then we don't need to actually make the action. Very often this is not the case.

We cannot predict exactly the future, but we can learn experience by actually executing actions in the environment and try to reconstruct the model.

Learning: given past experience $(z_{0:k})$, determine the model (f, h) .

Very important is the concept of state: a state can be seen as a snapshot of the dynamic system. The state can be either observable or not.

When the state is **fully observable**, the decision making problem for an agent is to decide which action must be executed in a given state.

The agent has to compute the function

$\pi: X \rightarrow A$ choose the best action to execute in a given state.
e.g. in a chess domain, given current configuration $x \in X$ choose the best move $a \in A$

When the model of the system is not known, the agent has to learn the function π .

How to learn the function?

■ Supervised learning

learning a function $f: X \rightarrow Y$ given a dataset $D = \{(x_i, y_i)\}$

■ Reinforcement learning

learning a behavior function $\pi: X \rightarrow A$ given a set of observations $D = \{(x_1, a_1, r_1, \dots, x_m, a_m, r_m)^{(i)}\}$

this dataset is not provided in terms of input/output pairs, we have instead the concept of rewards

A **reward** is something given to the agent during the evolution

representation of the evolution of the system

An action in supervised learning needs to be the best action

An action in reinforcement learning does not need to be the best action

a_m is not the best action we can do in x_m , but is instead the action we actually did in x_m .

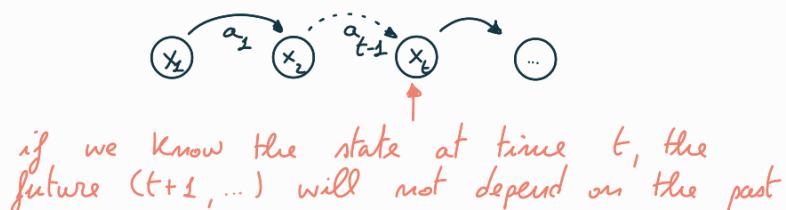
Representation of a dynamic system

X : set of states	$\left\{ \begin{array}{l} \text{explicit discrete and finite representation } X = \{x_1, x_2, \dots, x_n\} \\ \text{continuous representation } x = f(\dots) \text{ given by a state function} \\ \text{probabilistic representation } P(x) \text{ given by a probabilistic state function} \end{array} \right.$
A : set of actions	$\left\{ \begin{array}{l} \text{explicit discrete and finite representation } A = \{a_1, a_2, \dots, a_m\} \\ \text{continuous representation } A = u(\dots) \text{ given by a control function} \end{array} \right.$
δ : transition function	$\left\{ \begin{array}{l} \text{deterministic} \\ \text{non-deterministic} \\ \text{probabilistic} \end{array} \right.$
	<p>when an agent executes an action, there will be a probability distribution on what the outcome will be (successor state)</p>
Z : set of observations	$\left\{ \begin{array}{l} \text{explicit discrete and finite representation } Z = \{z_1, \dots, z_k\} \\ \text{continuous representation } z = z(\dots) \text{ given by an observation function} \end{array} \right.$
	<p>only used in Hidden Markov Models; here the state is fully observable</p>

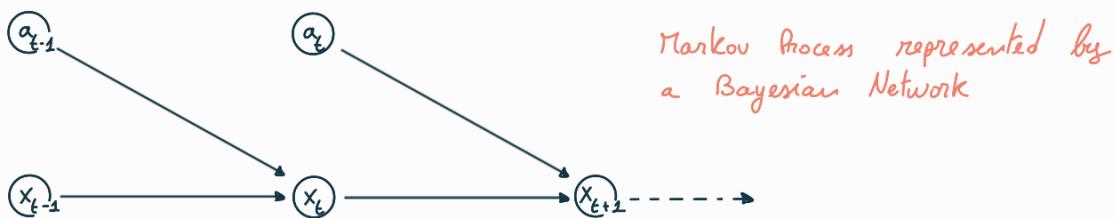
Markov Property

- Once the current state is known, the evolution of the dynamic system does not depend on the history of states, actions and observations;
- The current state contains all the information needed to predict the future;
- Future states are conditionally independent of past states and past observations given the current state;
- The knowledge about the current state makes past, present and future observations statistically independent

A Markov Process is a process that has the Markov property.



$$P(\text{future} | \text{current, past}) = P(\text{future} | \text{current})$$



x_{t+1} does not depend on x_{t-1} if we know x_t

Markov Decision Process (MDP)

MDP will model our dynamic system

$$MDP = \langle X, A, \delta, r \rangle$$

In a **deterministic** case:

- X is a finite set of states
- A is a finite set of actions
- $\delta : X \times A \rightarrow X$ is a transition function
- $r : X \times A \rightarrow R$ is a reward function

deterministic transition: any time we execute an action from a state we get the same result

In the **non-deterministic** case:

- X is a finite set of states
- A is a finite set of actions
- $\delta : X \times A \rightarrow 2^X$ → we get a set of possible states
- $r : X \times A \times X \rightarrow R$

different rewards: the reward also depends on the next state

In the **stochastic** case we have a probability distribution associated to the outcome of an action executed from a state



$$\delta = P(x_{t+1} | x_t, a_t)$$

Solving an MDP

Solving an MDP means finding a policy (function from states to actions)

$$\pi : X \rightarrow A$$

We want the **optimal policy**

For each state $x \in X$, $\pi(x) \in A$ is the optimal action to be executed in such state
Optimality criterion = maximize the **reward**.

A **reward** says how good the behaviour is for the agent.

- maximizing the **cumulative** reward gives the same importance to short term and long term rewards;
- maximizing the **cumulative discounted** reward is better; the reward collected in the future will be discounted with a discount factor $\gamma \in [0, 1]$; this means that the same reward is more valuable now than in the future

value of the policy starting from state x_1 $\text{V}^\pi(x_1) = E[\gamma r_1 + \gamma^2 r_2 + \gamma^3 r_3 + \dots]$

Expected value for non-deterministic MDPs

This enforces the agent to gain rewards as soon as possible

$V^{\pi_1}(x_1) > V^{\pi_2}(x_2)$ means that policy π_1 is better than π_2 since it has a higher cumulative reward. The optimal policy is the best one, the one with higher cumulative discounted value.

Optimal policy: $\pi^* = \operatorname{argmax}_{\pi} V^{\pi}(x), \forall x \in X$

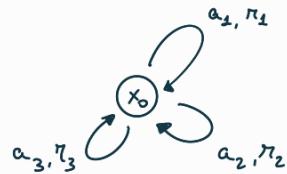
There may exist several optimal policies.

V^* is the value of the optimal policy.

One state MDP

A one-state MDP only has one state: x_0

- x_0
- A is a finite set of actions
- $\delta(x_0, a_i) = x_0 \quad \forall a_i \in A$ transition function
- $r(x_0, a_i, x_0) = r(a_i)$ reward function



Optimal policy: $\pi^*(x_0) = a_i$

This is also called **the bandit problem**.

This problem is interesting since its solution is a skeleton for most RL problems:

- ① Initialize a data structure \emptyset

$\theta_{(0)}[i] \leftarrow 0$ and $c[i] \leftarrow 0, i = 1 \dots |A|$

↓ action index
↓ time index

- ② for each time $t = 1, \dots, T$ (until termination condition)

- ① choose an index \hat{i} for action $a_{(t)} = a_{\hat{i}} \in A$
- ② execute action $a_{(t)}$ and collect reward $r_{(t)}$
- ③ increment $c[\hat{i}]$
- ④ update $\theta_{(t)}[\hat{i}] \leftarrow \frac{1}{c[\hat{i}]} (r_{(t)} + (c[\hat{i}] - 1) \theta_{(t-1)}[\hat{i}])$

- ③ optimal policy: $\pi^*(x_0) = a_{\hat{i}}$, with $\hat{i} = \operatorname{argmax}_{i=1 \dots |A|} \theta_{(T)}[i]$

How to find the optimal policy?

- Value iteration (estimate the value function and then compute π)
- Policy iteration (estimate directly π)

Learning through value iteration

The agent could learn the value function $V^*(x)$
From V^* it can determine the optimal policy:

the optimal policy is the policy that maximizes this quantity

$$\pi^*(x) = \arg\max_{a \in A} [r(x, a) + \gamma V^*(\delta(x, a))]$$

$\underbrace{r(x, a) + \gamma V^*(\delta(x, a))}_{Q(x, a)}$

However this policy cannot be computed this way because δ and r are not known.
Note that:

$$V^*(x) = \max_{a \in A} \{ r(x, a) + \gamma V^*(\delta(x, a)) \} = \max_{a \in A} Q(x, a)$$

thus we can rewrite

$$Q(x, a) = r(x, a) + \gamma V^*(\delta(x, a))$$

transition function

as:

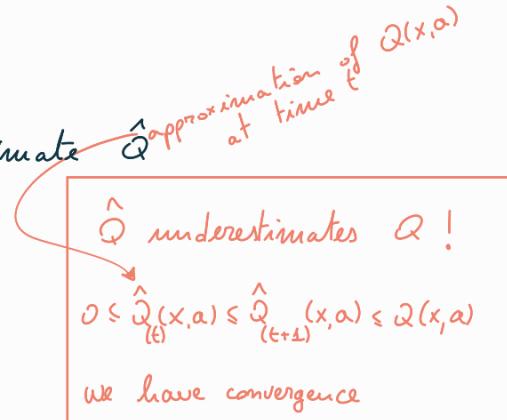
$$Q(x, a) = r(x, a) + \gamma \max_{a'} (Q(\delta(x, a), a'))$$

instead of using the real Q function, we use an estimate

$$\hat{Q}(x, a) \leftarrow \bar{r} + \gamma \max_{a'} \hat{Q}(x', a')$$

↓ immediate reward we get when we execute a from x

estimate



At the next iteration we will have an estimate $\hat{Q}(x, a)$
We don't know the reward function until we execute the action

Algorithm:

- ① for each x, a initialize table entry $\hat{Q}_{(0)}(x, a) \leftarrow 0$
- ② observe current state x
- ③ for each time $t = 1, \dots, T$ (until termination condition)
 - choose an action a
 - execute the action a
 - observe the new state x'
 - collect immediate reward \bar{r}
 - update the table entry for $\hat{Q}(x, a)$ as follows:

$$\hat{Q}_{(t)}(x, a) \leftarrow \bar{r} + \gamma \max_{a' \in A} \hat{Q}_{(t-1)}(x', a')$$

■ $x \leftarrow x'$

- ④ Optimal policy: $\pi^*(x) = \arg\max_{a \in A} \hat{Q}_{(T)}(x, a)$

Note: not using δ and r , but just observing new state x' and immediate reward \bar{r} after the execution of the chosen action

How actions are chosen by the agents?

- Exploitation : select action a that maximizes $\hat{Q}(x, a)$
- Exploration : select random action a (with low value of $\hat{Q}(x, a)$)

ϵ -greedy strategy

Given $0 \leq \epsilon \leq 1$

- select a random action with probability ϵ
- select the best action with probability $1-\epsilon$

ϵ can decrease over time (first exploration, then exploitation)

Soft-max strategy

Actions with higher \hat{Q} values are assigned higher probabilities, but every action is assigned a non-zero probability

↳ in order to have convergence

$$p(a_i|x) = \frac{\hat{Q}(x, a_i)}{\sum_j K^{\hat{Q}(x, a_j)}}$$

$K > 0$ determines how strongly the selection favors actions with high \hat{Q} values.

K may increase over time (first exploration, then exploitation).